

# Data Wrangling

What is data wrangling?

It is the process of cleaning, organizing & transforming raw data into usable format for analysis.

- creating new columns
- handling large amount of messy data

Need to load

→ Tidyverse

- install packages
- 8 diff packages

>Select

→ select column (mostly use)

: → between columns to select all column (in order)

code:

select() function

Renome <- select(data, column1, column2 : column5,  
column 6)

## Filter function [filter( )]

→ Subsetting / filtering data

- subset data fairly easy

Example

filtering Treatment "conv" from data

head(filter(data, Treatment == "conv"))

↓  
head will display first 6 lines

complex using %>% or |

head(filter(data, Treatment == "conv") |

↓  
filter(Fungicide == "c"))

Display first 6 lines

\* function mutate( )

- To create new column

Code converting thickness to log-transform

Base R

data\$Newcolumn <- log(data\$thickness)

Advanced

head(mutate(data, Newcolumn = log(thickness)))

## ④ Pipe

→ Combining multiple functions together

→ combine output from one function  
into input of another function

→ Input the previous data from the  
right side of the pipe into the  
function on the left %>%

Example:

combining whatever we've done so far

(mutate, filter, select)

↑ facing right  
data %>% select (

- input this data frame to this function
- don't need to supply data frame  
in each function using pipes

what it's doing final code

filtering ← data %>% select( column1, column2, column3:  
column5, column6 ) %>%

subsetting ← filter( column4 = "conv" ) %>%

Creating new column ← mutate ( Newcolumn = log( Oldcolumn ) ) %>%

head ( )

## Summarise data

→ To find sd, mean etc.

Code:

```
data %>% select ( ) %>%
filter ( ) %>%
mutate ( ) %>%
summarise (outputname = mean (New column))
```

Output

Mean.Rich / output name

| → mean of the new column

# connecting multiple summarise statistics

Previous pipeline %>%

summarise (Mean.Rich = mean (New column),  
↓  
output nm)

n = n ( ),

sd.dev = sd (New column) %>%

To create ← mutate (std.error = sd.dev / sqrt(n))  
a new  
column of  
sd.error.

Output

Mean.Rich	n	sd.dev	std.error
X	y		Z

function: group-by()

Summarise data by group

→ Calculating summary statistics by group

Code:

```
data %>%  
  select( ) %>%  
  group_by( var1, var2) %>%  
  mutate( Newcol = log(old col)) %>%  
  summarise( ) ;  
  n = n()  
  sd.dev = sd( Newcol)  
  mutate( std.error = sd.dev/sqrt(n))
```

connecting to plotting

→ output of plot and connect to ggplot

Previous code %>%

```
ggplot( aes(x= var1, y= var2))
```

```
geom_bar( stat = "var1")
```

```
geom_errorbar( aes(x= var1, ymin= var2, st.dev  
+ forax = var2 + std.error), width=0.4)
```

```
+ theme_minimal()
```

```
xlab(" ")
```

```
ylab("Label") +
```

```
facet_wrap(~ var4)
```

## Joining function

→ `left-join()`

Keep all rows of X, add matching from Y

→ `right-join()`

→ `right-join` ⇒ all Y, match X

(reverse of `leftjoin`)

→ `inner-join`

only keeps rows common to both X and Y

→ `full-join`

→ keep all from X & Y

→ combining various set of data (?)  
based on similar column

code:

```
# selecting just richness and sample ID
```

```
richness <- data %>%
```

```
select (sample ID, richness)
```

```
# selecting column that don't include richness
```

```
metadata <- data %>%
```

```
select (Sample ID, other columns)
```

```
# Adding the richness data to the metadata based on  
# A the common column of sample ID
```

```
head (left_join (metadata, richness, by = "SampleID"))
```

## Pivot function

→ converting from wide to long data (one observation per sample)

i) `-pivot-wide`

converts long format to wide

ii) `pivot-long`

wide → long

What is wide or long data

Long data

- Each row rep. a single observation

- values stack in one column

- for statistical analysis

Wide data

- Each row represents single subject or experimental unit

- Diff categories are stored in separate column.